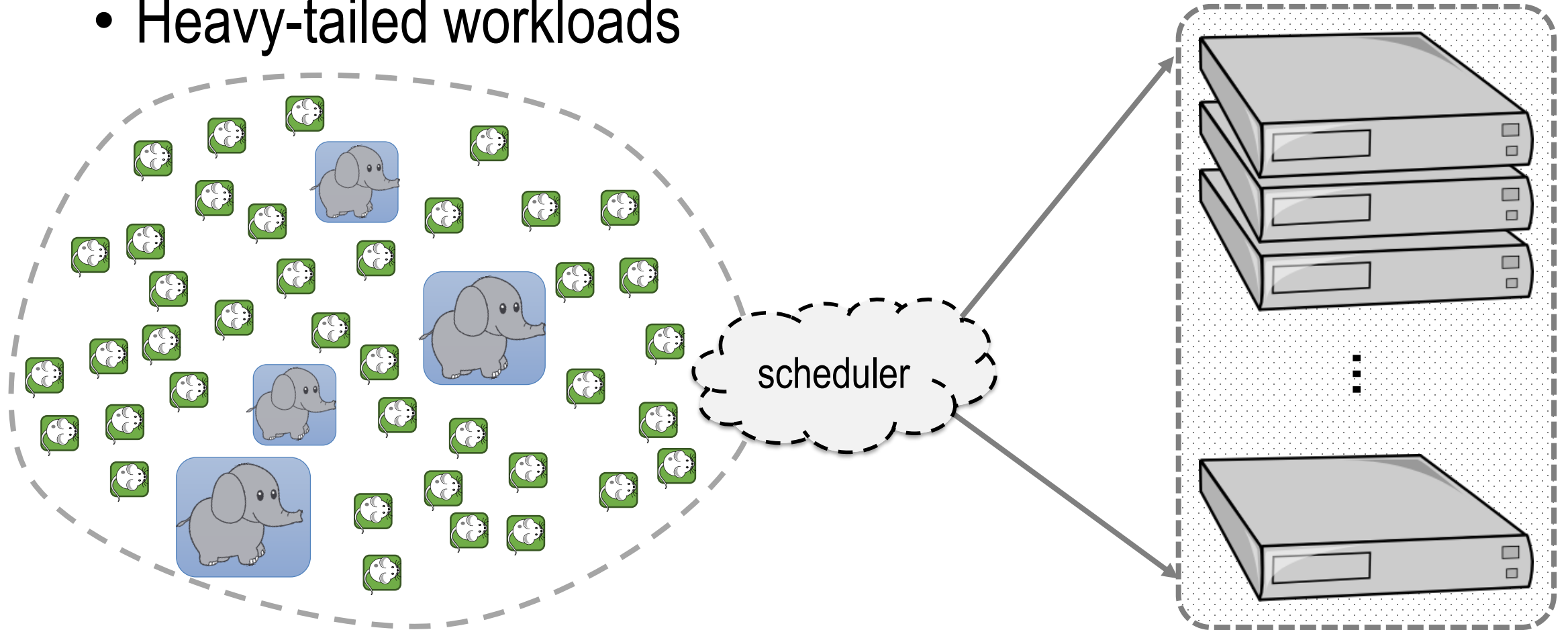# Kairos

**Data center scheduling
without task runtime estimates**

# Kairos key idea

- New preemption approach


✓ No head-of-line blocking

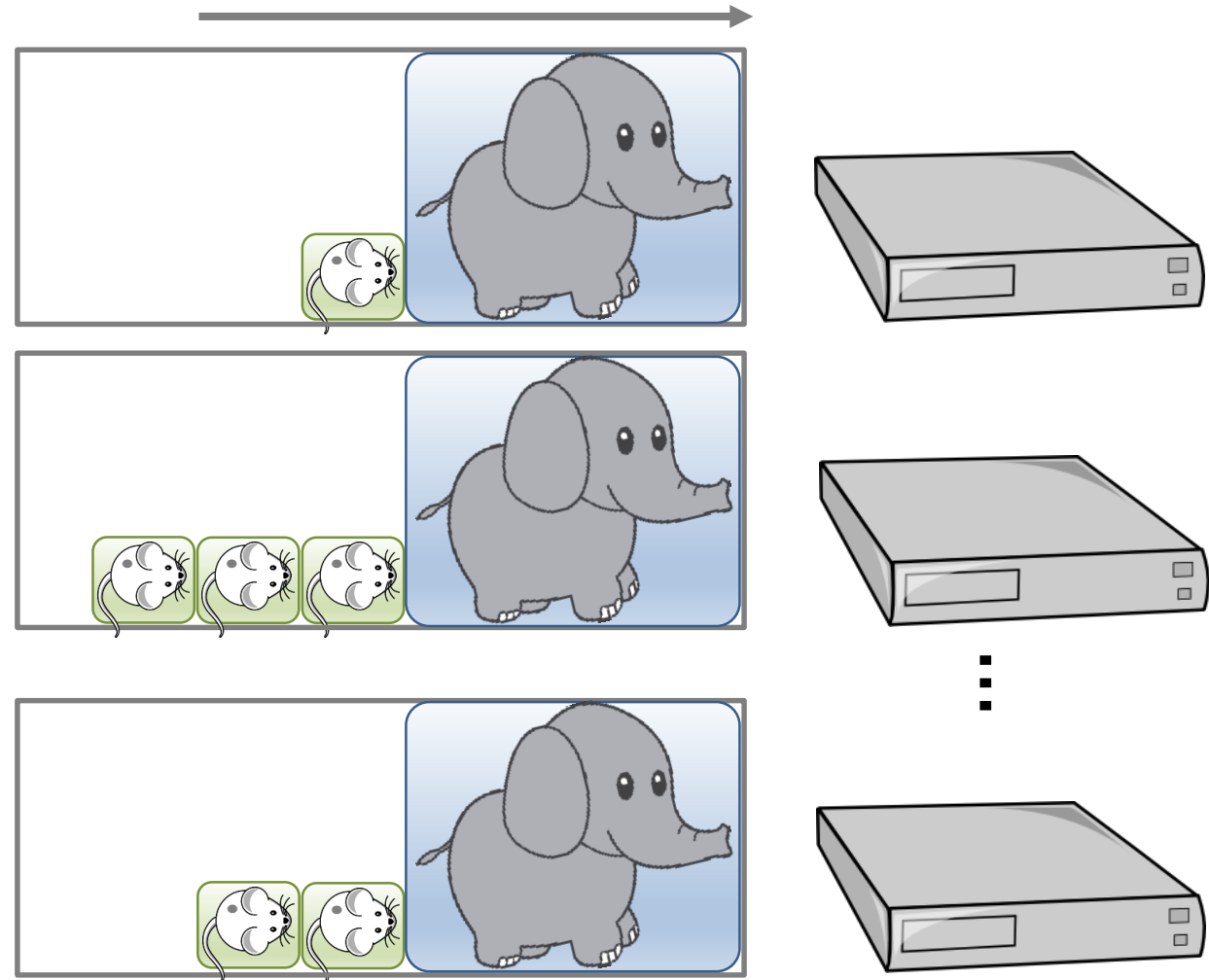✓ Good scheduling performance

# Data center scheduling challenge

- Heavy-tailed workloads
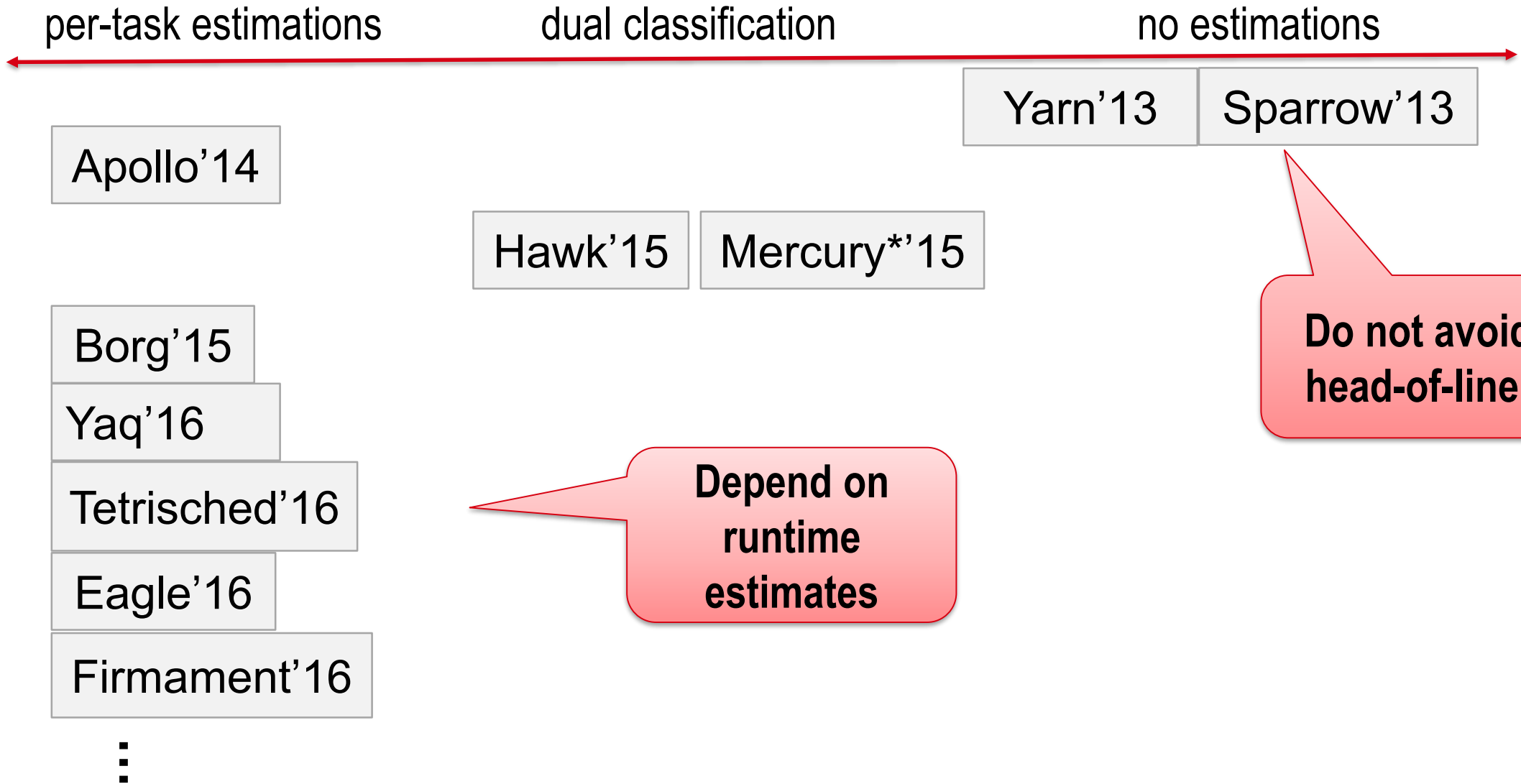
cluster



scheduler

Kairos: Preemptive Data Center Scheduling Without Runtime Estimates | SoCC'18    4

# Problem: head-of-line blocking

- Short waiting for long
- High likelihood

# Historical use of runtime estimates

per-task estimations          dual classification          no estimations

Yarn'13 | Sparrow'13

Apollo'14

Hawk'15 | Mercury*'15

**Do not avoid head-of-line!**

Borg'15

Yaq'16

Tetrisched'16

**Depend on runtime estimates**

Eagle'16

Firmament'16

⋮

# Hard to obtain reliable estimates

- Mis-estimations happen
  - unseen jobs, skewed input, failures/spikes

- Consequences:
  - poor scheduling decisions*, violate SLOs^
  - complex designs to compensate

*Job-aware scheduling in Eagle: Divide and Stick to Your Probes (SoCC'16)
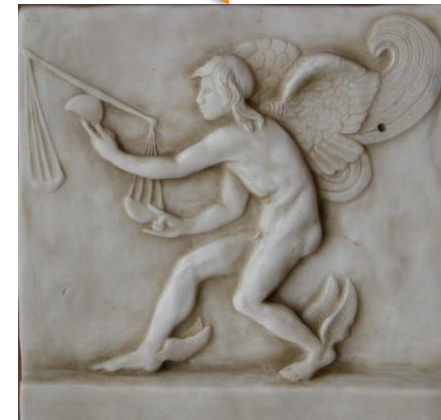^ Tetrisched: global rescheduling with adaptive plan-ahead in dynamic heterogeneous clusters (Eurosys'16)

Can we dispense with task runtime estimates altogether?

Can we dispense with task runtime estimates altogether?
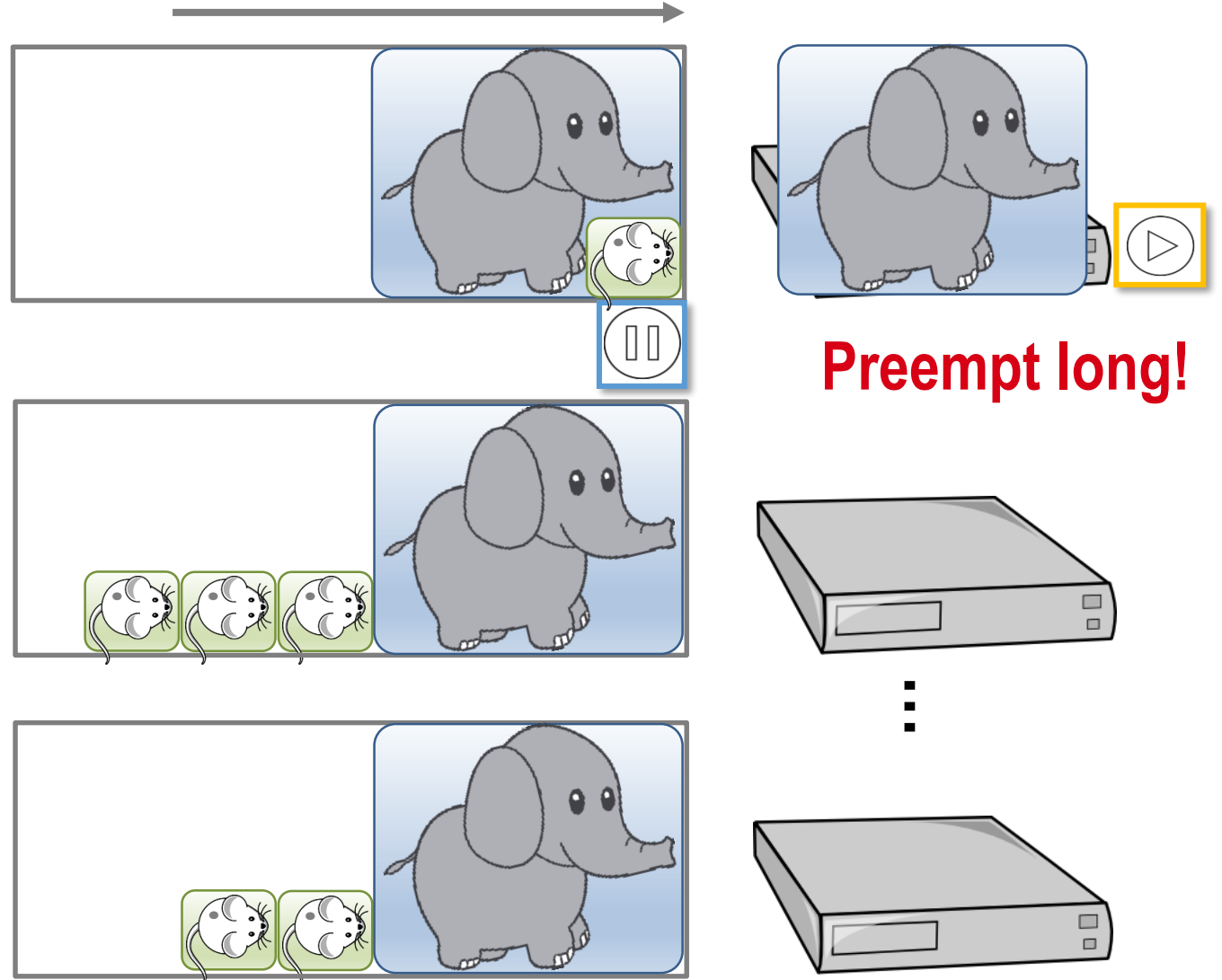
✓ Avoid head-of-line blocking
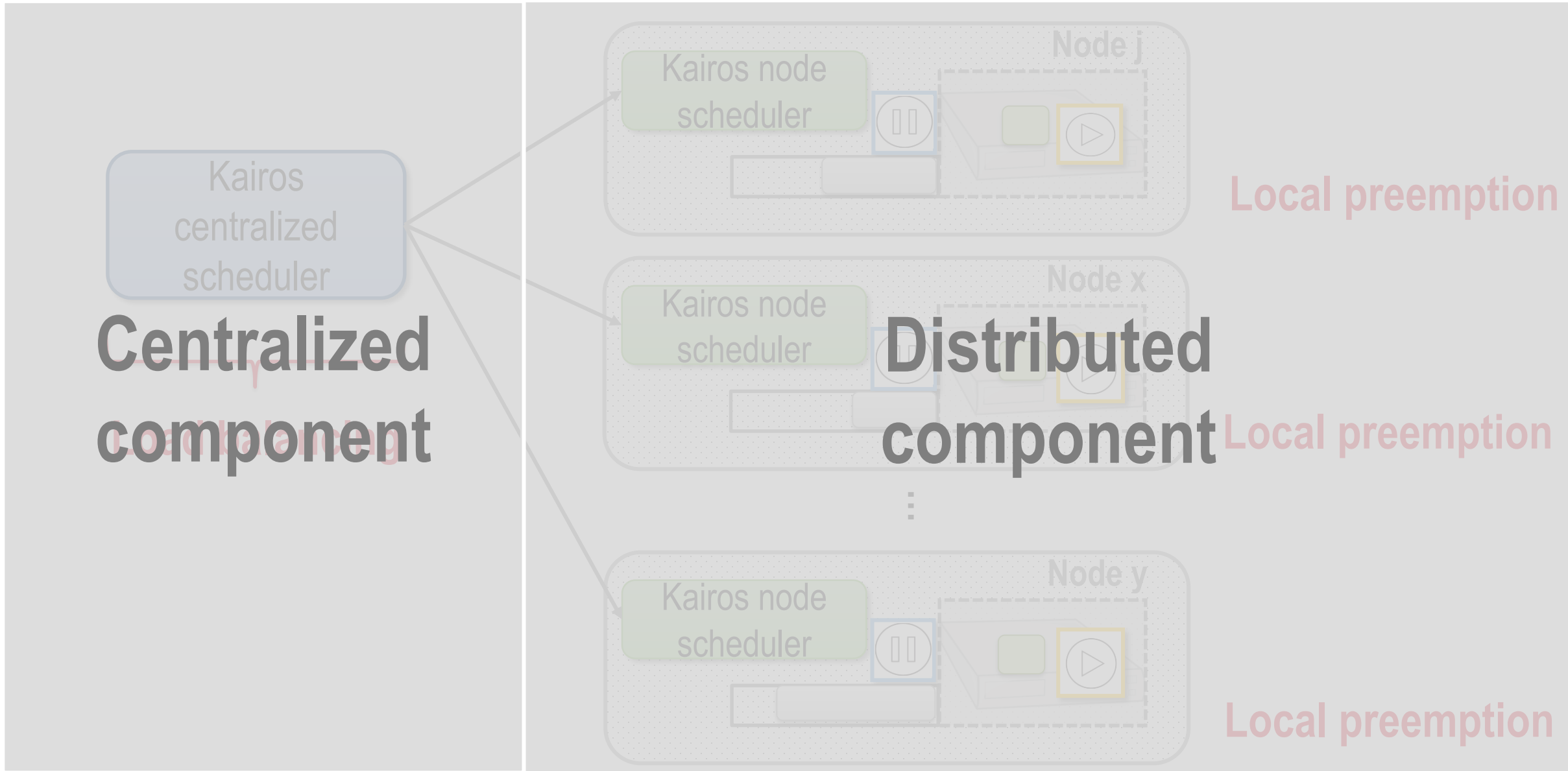✓ No task runtime estimates
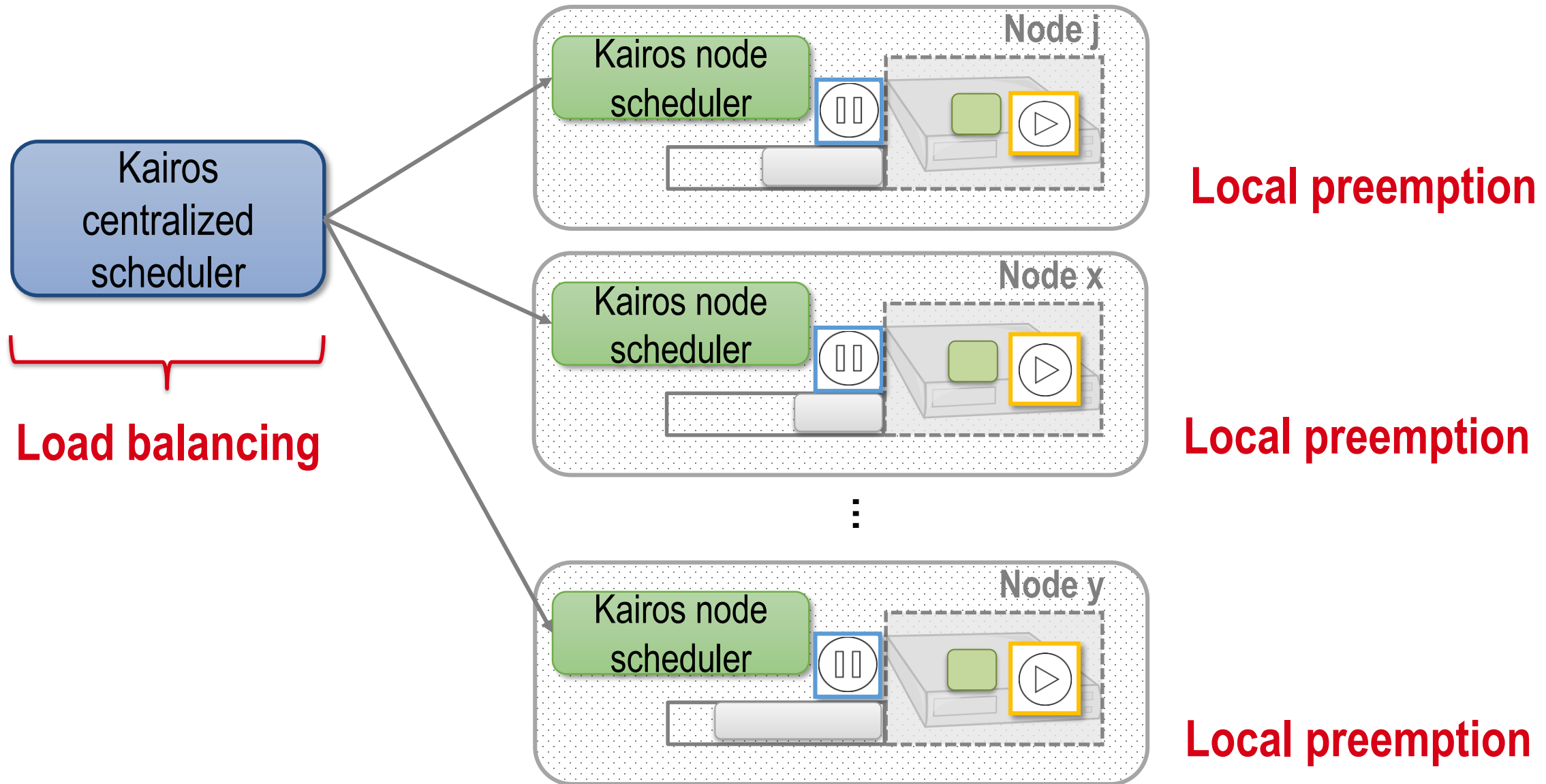
**Kairos**

# Kairos insight

Use preemption!!

# Preemption in Kairos

Costly resuming elsewhere:
Do preemption locally!



**Preempt long!**

# Kairos architecture



Kairos centralized scheduler

**Centralized component**

Kairos node scheduler — Node j

Local preemption

Kairos node scheduler — Node x

**Distributed component**

Local preemption

Kairos node scheduler — Node y

Local preemption

# Kairos architecture



Kairos centralized scheduler

Kairos node scheduler — Node j

**Local preemption**

Kairos node scheduler — Node x

**Local preemption**

Kairos node scheduler — Node y

**Local preemption**

**Load balancing**

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Kairos architecture



Kairos centralized scheduler

Load balancing

Node j

Kairos node scheduler

**Local preemption**

Node x

Kairos node scheduler

**Local preemption**

Node y

Kairos node scheduler

**Local preemption**

# Least-Attained Service (LAS)

- Preemptive policy

- Give resources to task that received least service

- ✓ New task runs immediately

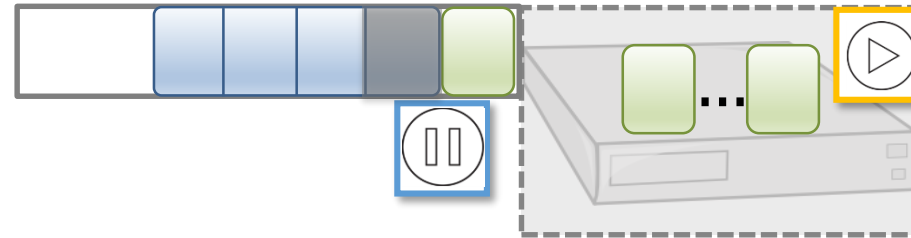- ✓ Runs as long as it is the one with least received service

# LAS rationale

- Good for heavy-tailed workloads*

- Benefits:

1. Shorter tasks have priority (no head-of-line blocking)
2. Shorter tasks –very likely– execute until completion

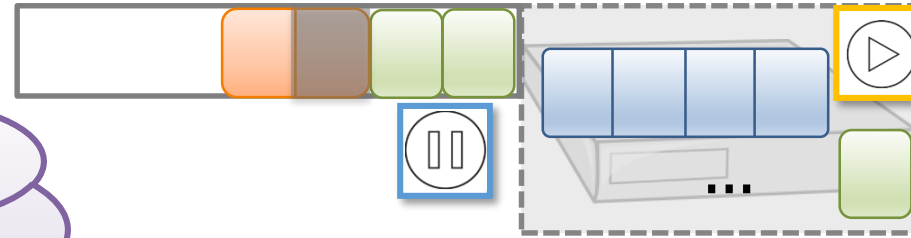*Performance modeling and design of computer systems: queueing theory in action M. Harchol-Balter 2013

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Kairos distributed scheduling

- Node schedulers
  - LAS at the nodes

How to dispatch tasks among nodes?
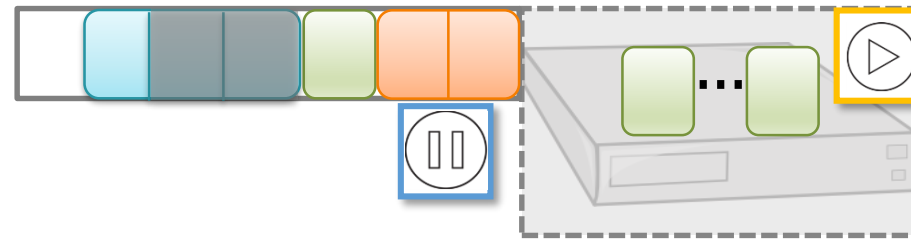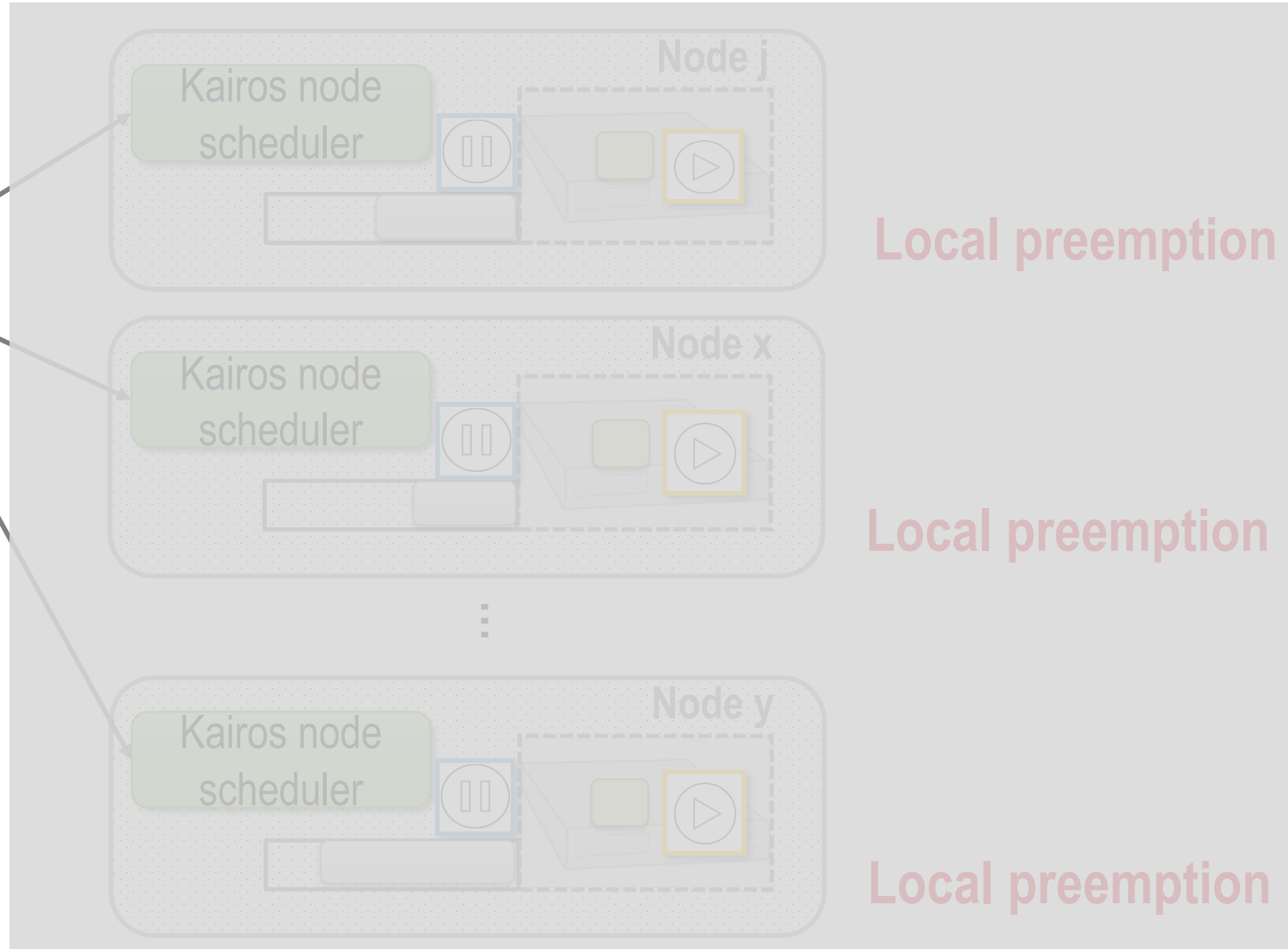
Kairos node scheduler

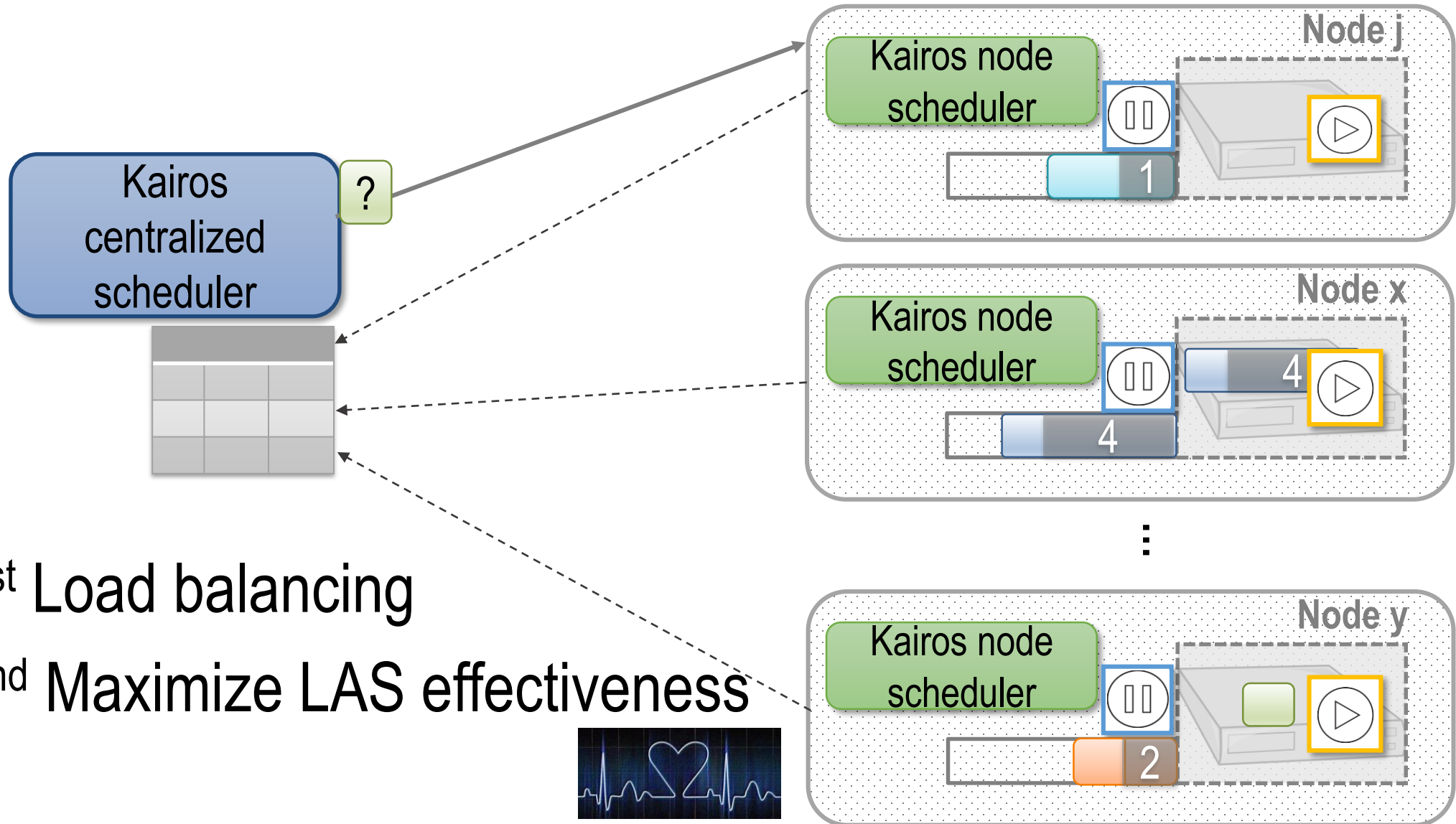Kairos node scheduler

Kairos node scheduler

# Kairos architecture

Kairos centralized scheduler

**Load balancing**

Node j

Kairos node scheduler

Local preemption

Node x

Kairos node scheduler

Local preemption

Node y

Kairos node scheduler

Local preemption

ÉCOLE POLYTECHNIQUE
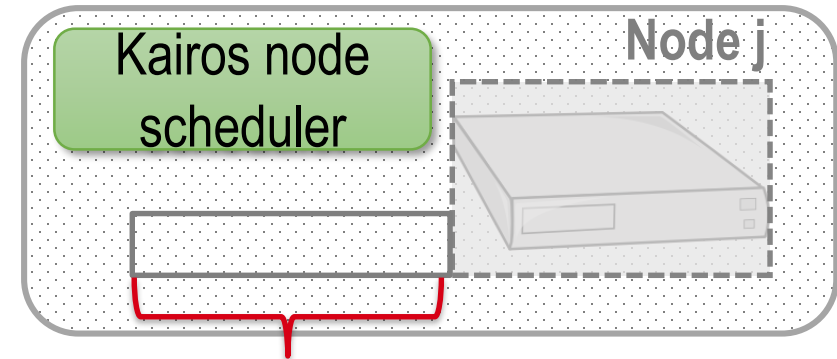FÉDÉRALE DE LAUSANNE

# Kairos centralized scheduling



1st Load balancing

2nd Maximize LAS effectiveness

# Load balancing rationale

1. Lowest # tasks: no idle nodes
- Bound max # tasks

**1. Avoid!**

Node j

Kairos node scheduler

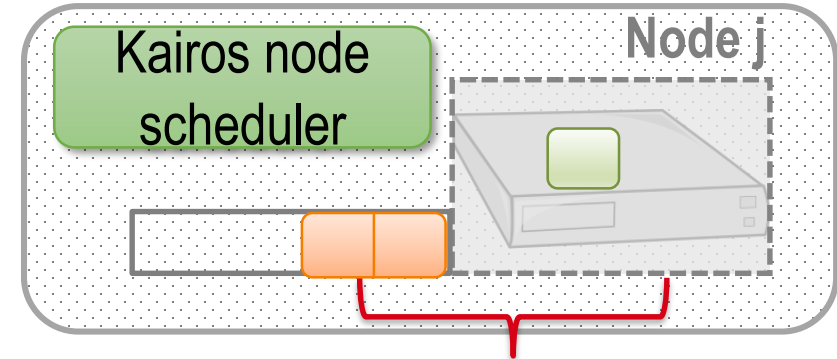**0 tasks**

Node y

Kairos node scheduler
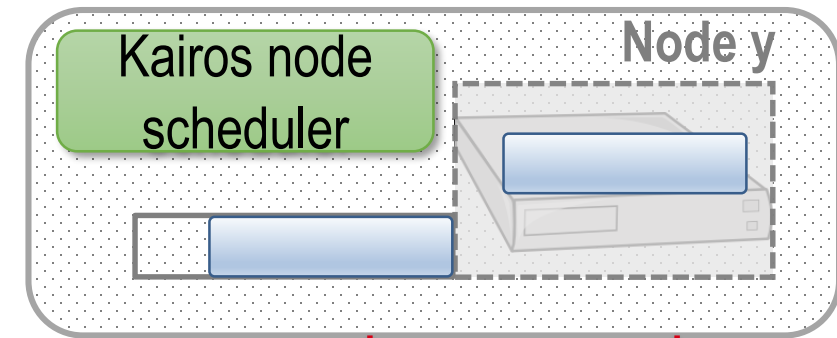
**100 tasks**

# Load balancing rationale

2. LAS-aware policy break ties:

- Heavy-tailed for each node

- Maximize LAS effectiveness

- Node with lowest AS variance*

**2. Avoid!**



only short

only long

*Minimizing total flow time and total completion time with immediate dispatching. Avrahami et.al. 2003*
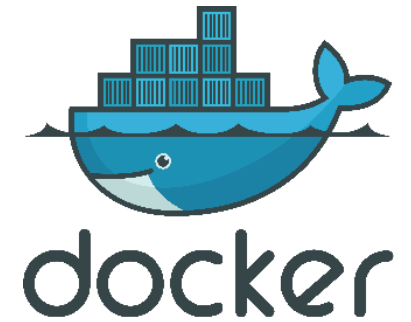
*Multi-layered round robin routing for parallel servers Down et.al. 2006*

# Kairos recap

1. Distributed:
   - ✓ LAS node level

2. Centralized:
   - ✓ LAS-aware load balancing technique

# Evaluation

- Yarn and Docker containers

- 120 cores in 30 nodes

- heavy-tailed workload (100 jobs)

- Metrics: Job runtime and slowdown

- Compare to: Big-C [ATC'17], FIFO

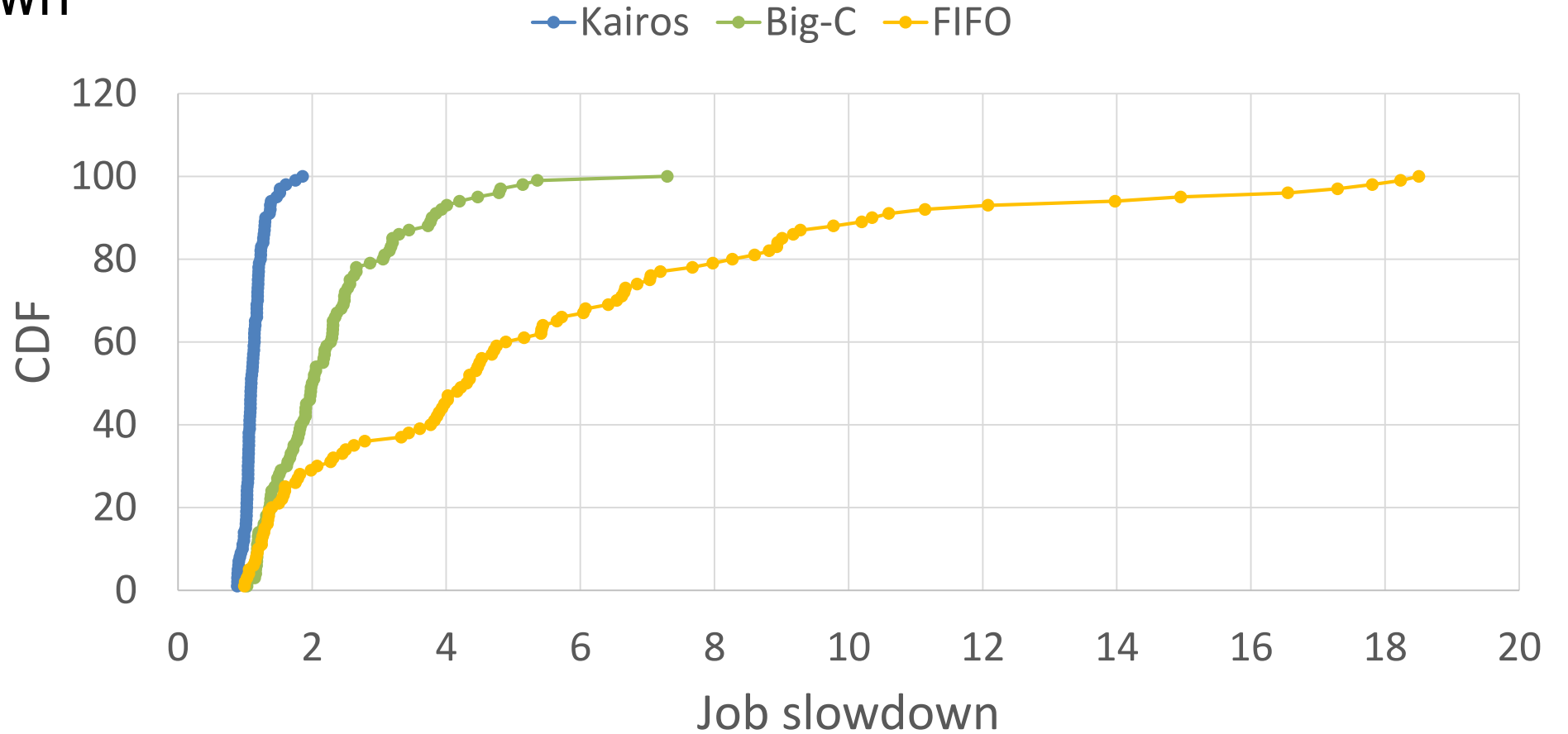- Simulation: Google trace, compare to Eagle [SoCC'16]

# What is the slowdown?

$$job\ slowdown = \frac{observed\ job\ runtime}{uncontended\ job\ runtime}$$

## Best job slowdown = 1

# Kairos vs Big-C and FIFO

Job slowdown
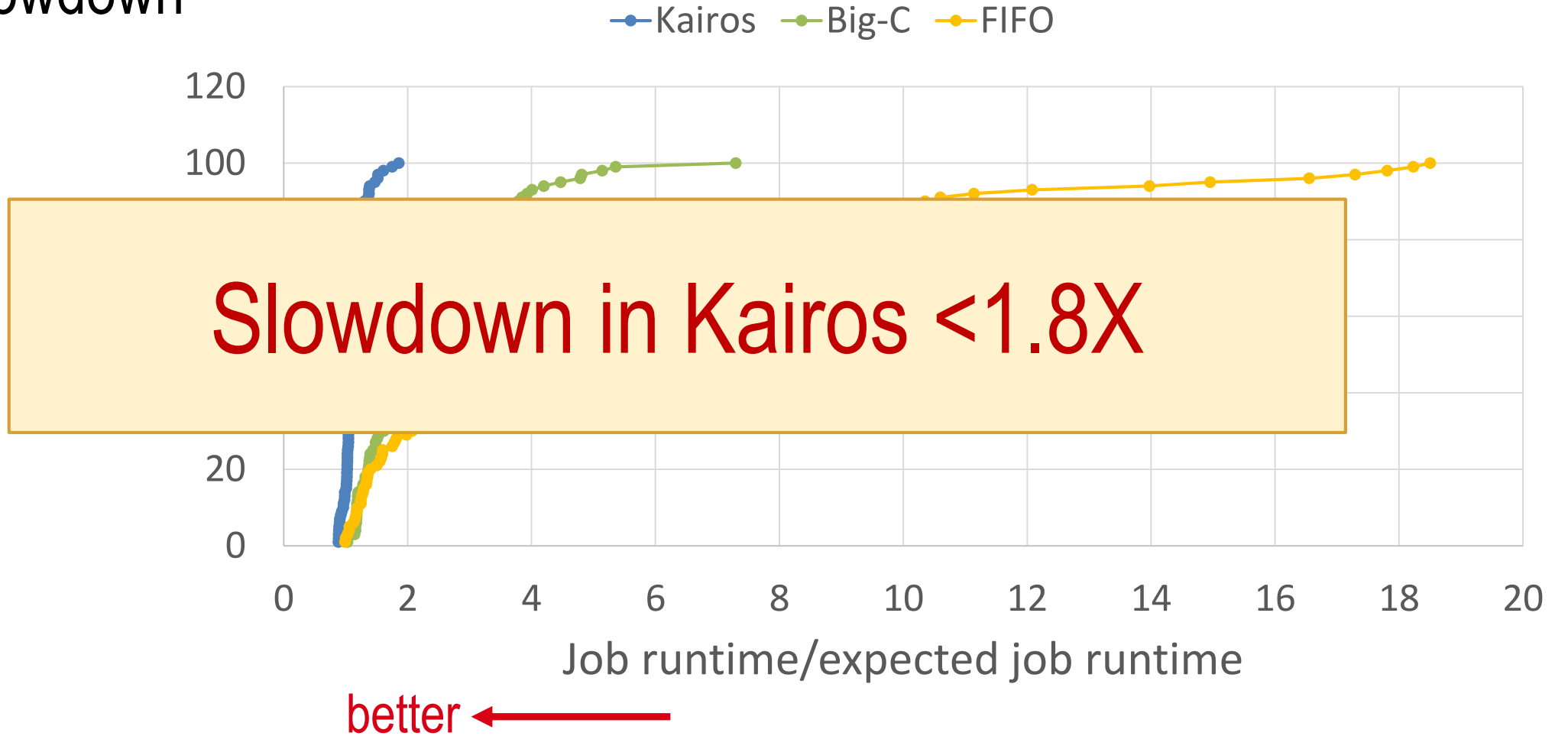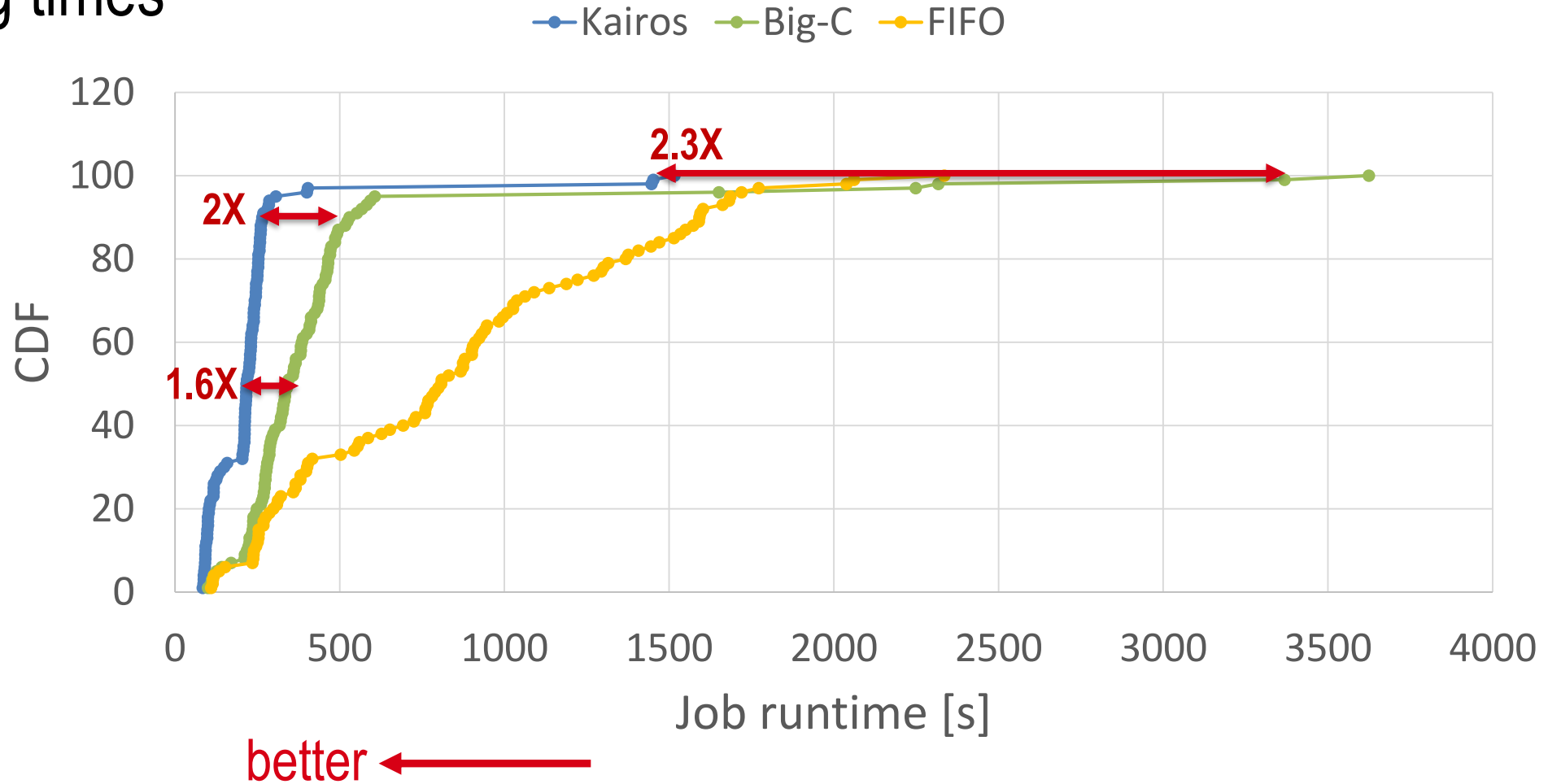


better ←

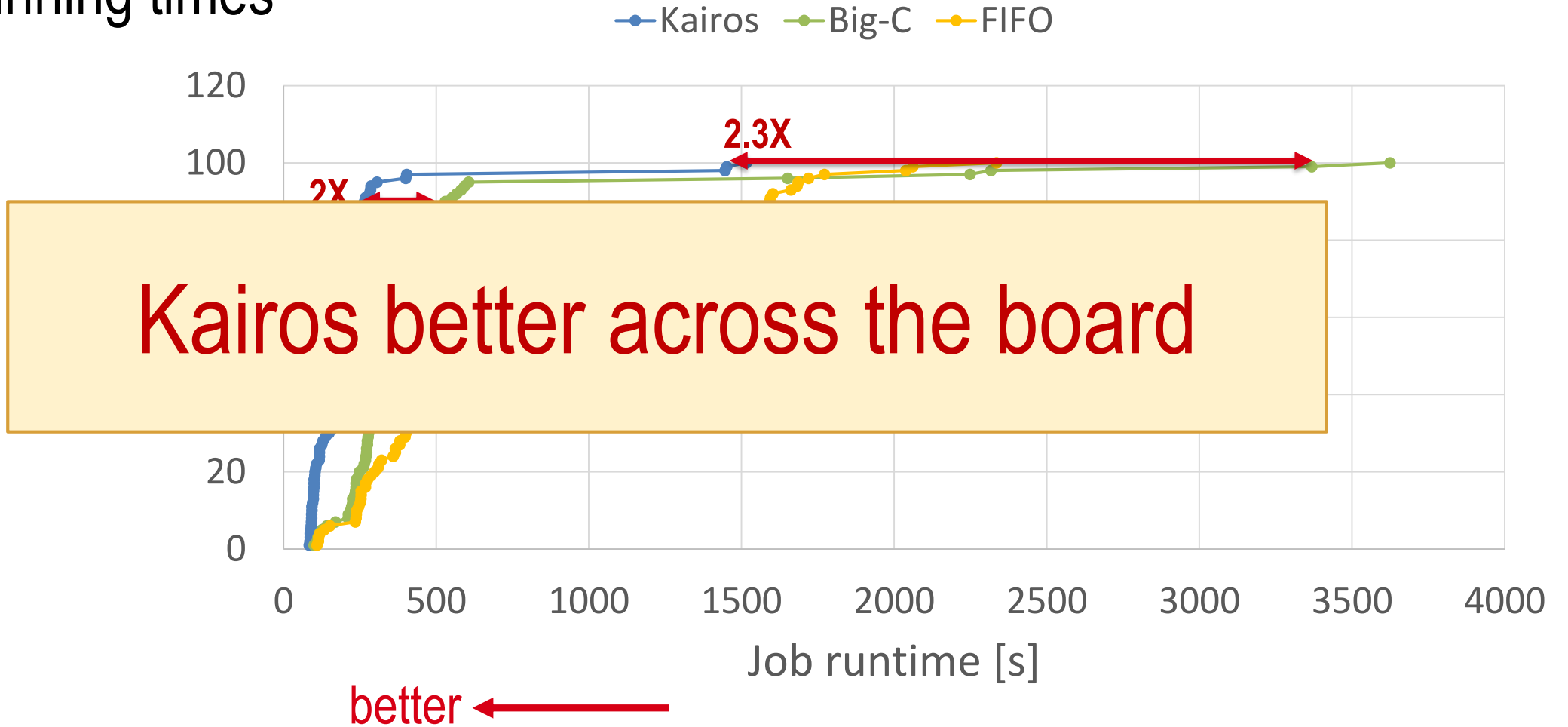# Kairos vs Big-C and FIFO

Job slowdown



Slowdown in Kairos <1.8X

Job runtime/expected job runtime

better ←

# Kairos vs Big-C and FIFO

Job running times

# Kairos vs Big-C and FIFO

Job running times



Legend: Kairos, Big-C, FIFO

2.3X

2X

**Kairos better across the board**

Job runtime [s]

better

# Kairos vs Eagle

- Short jobs runtime
- Google trace

# Kairos vs Eagle

- Short jobs runtime
- Google trace



**Kairos works well at large scale**

50th    90th    99th

1,4

1,2

0,4

0,2

0

better

Lower    Cluster load    Higher

# Why are we better?

Against FIFO

✓ FIFO does not avoid head-of-line

Against Big-C

✓ We do preemption better

Against Eagle

✓ Preemption

# More in the paper

- Evaluation with a uniform workload

- Sensitivity to parameters

- Comparison with other load balancing techniques

- How we do preemption

- Soon open sourced

# Kairos

✓ First preemptive scheduler without runtime estimates

✓ Smart preemption: good job runtime and slowdown

✓ LAS at node level
✓ LAS-aware load balancing